

Data-Driven Documentation

Multilingual Technology for Producers of Information

Aarne Ranta
Digital Grammars AB

12 April 2016

digital  grammars
Language technology to rely on.

Problem:

reliable and efficient translation

Machine translation is sometimes good, sometimes bad - and you never know how it will be this time.

English Swedish Finnish Detect language ▼



Dutch Chinese (Simplified) English ▼

Translate

Min mor är inte svensk.

我的母亲是瑞典的。

English Swedish Finnish Detect language ▾



Dutch Chinese (Simplified) English ▾

Translate

Min mor är svensk.
Min mor är inte svensk.



我的母亲是瑞典的。
我的母亲是瑞典的。

English Swedish Finnish Detect language ▾

Min mor är svensk.
Min mor är inte svensk.



Dutch Chinese (Simplified) English ▾

Translate

我的母亲是瑞典的。
我的母亲是瑞典的。

English Swedish Finnish Detect language ▾

Min mor är svensk.
Min mor är inte svensk.



Dutch Arabic English ▾

Translate

My mother is Swedish.
My mother is Swedish.

English Swedish Finnish Detect language ▾

Min mor är svensk.
Min mor är inte svensk.



Dutch Chinese (Simplified) English ▾

Translate

我的母亲是瑞典的。
我的母亲是瑞典的。

English Swedish Finnish Detect language ▾

Min mor är svensk.
Min mor är inte svensk.



Dutch Arabic English ▾

Translate

My mother is Swedish.
My mother is Swedish.

Min far är svensk.
Min far är inte svensk.



My father is Swedish.
My father is not Swedish.

Who cares?

Consumer translator:

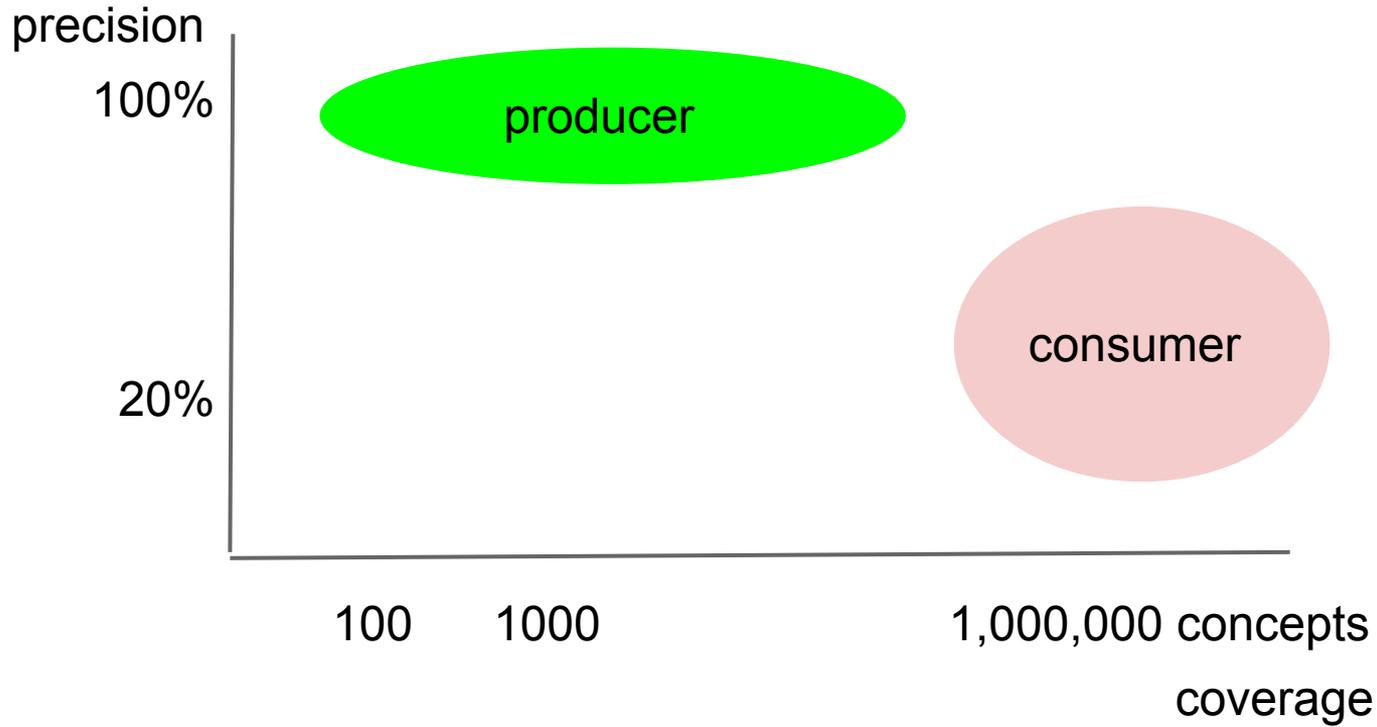
- browsing quality: to get an idea
- reader is responsible
- + translate anything

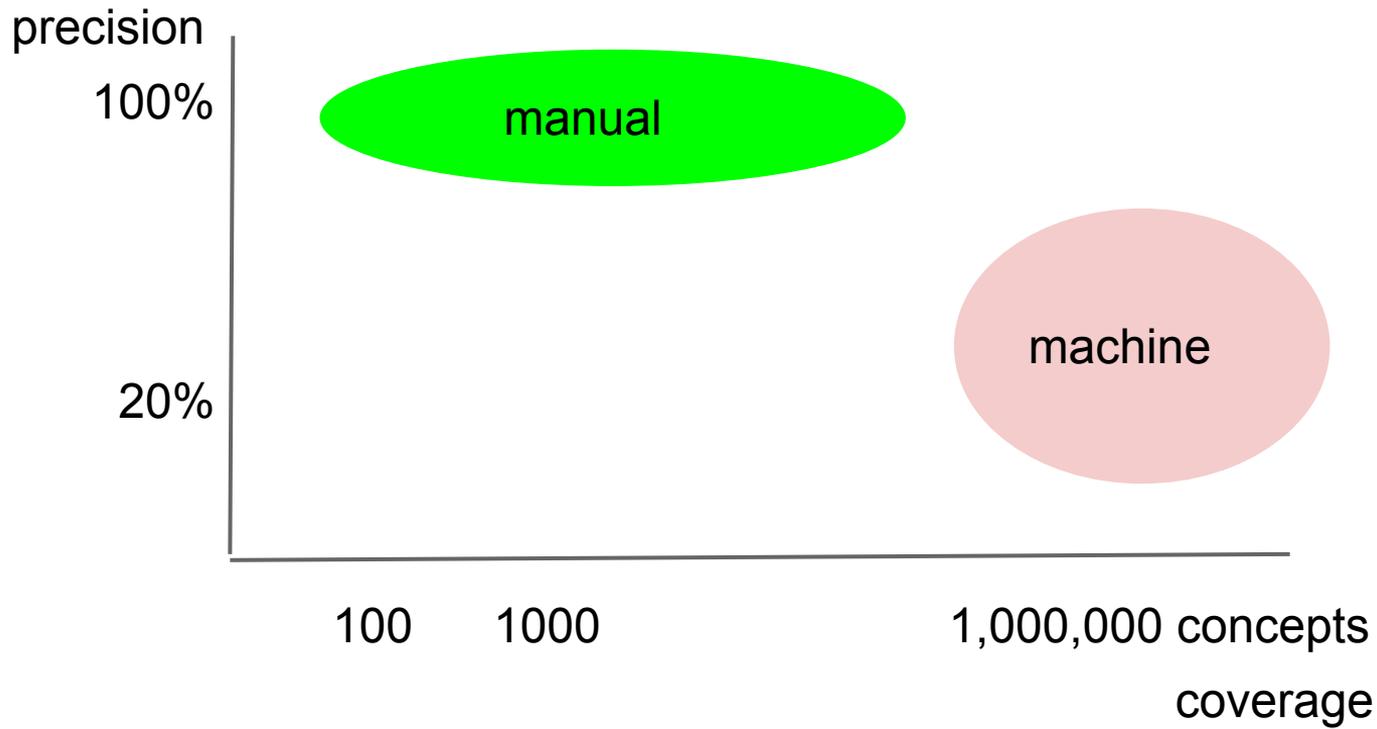
Consumer translator:

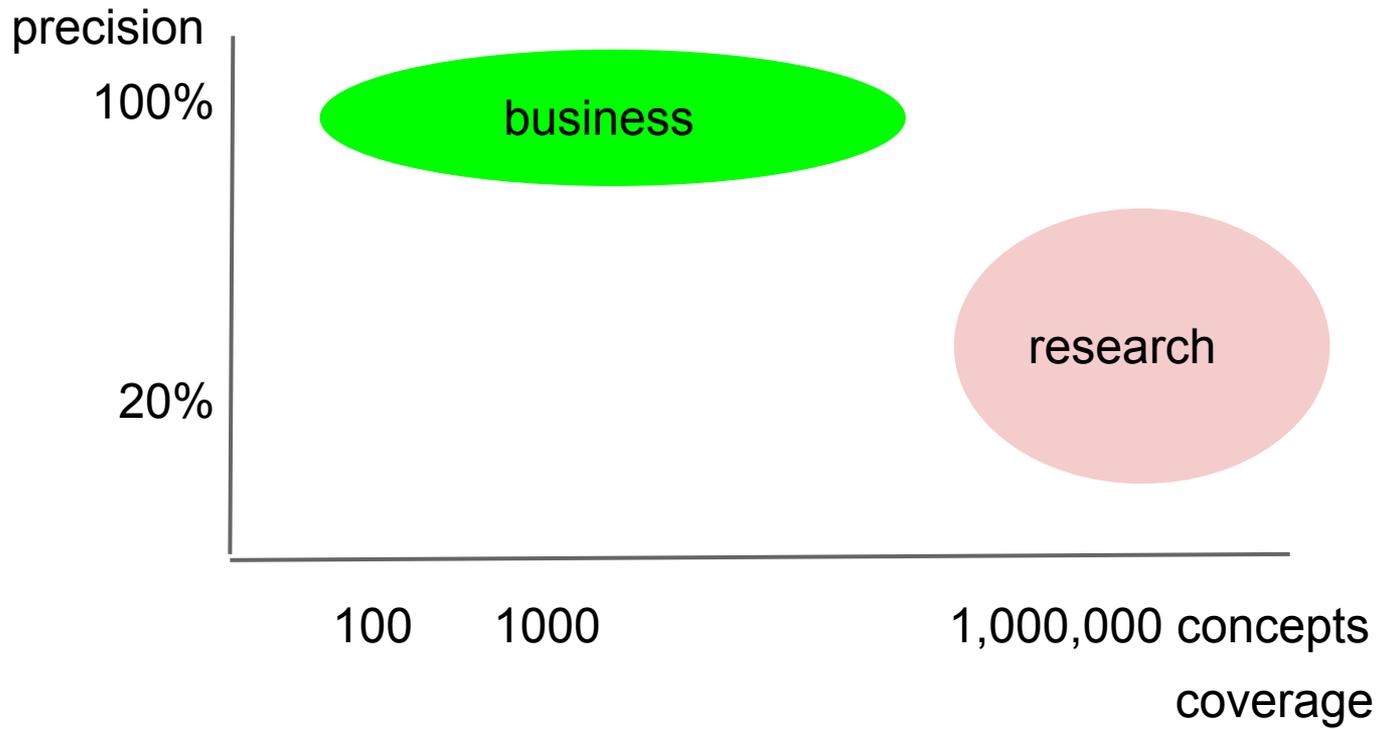
- browsing quality: to get an idea
- reader is responsible
- + translate anything

Producer translator:

- + publication quality: to get everything right
- + publisher is responsible
- translate my content







A solution:

Data-Driven Documentation

digital Grammars

Language Technology

2014 -

ely on.

REMU

VR 2013 - 2017

CLT

2009 - 2015

MOLTO

EU 2010 - 2013

G

1998 -

Data

object	property	value
door	free width	121cm
walking area	tilt sideways	0.5%

Data

object	property	value
door	free width	121cm
walking area	tilt sideways	0.5%

Documentation: Eng

The free width of the door is 121cm.

The walking area tilts 0.5% sideways.

Data

object	property	value
door	free width	121cm
walking area	tilt sideways	0.5%

Documentation: Eng

The free width of the door is 121cm.
The walking area tilts 0.5% sideways.

Documentation: Swe

Dörrens fria bredd är 121cm.
Gångytan lutar 0.5% i sidled.

Data

object	property	value
door	free width	121cm
walking area	tilt sideways	0.5%

Documentation: Eng

The free width of the door is 121cm.
The walking area tilts 0.5% sideways.

Documentation: Swe

Dörrens fria bredd är 121cm.
Gångytan lutar 0.5% i sidled.

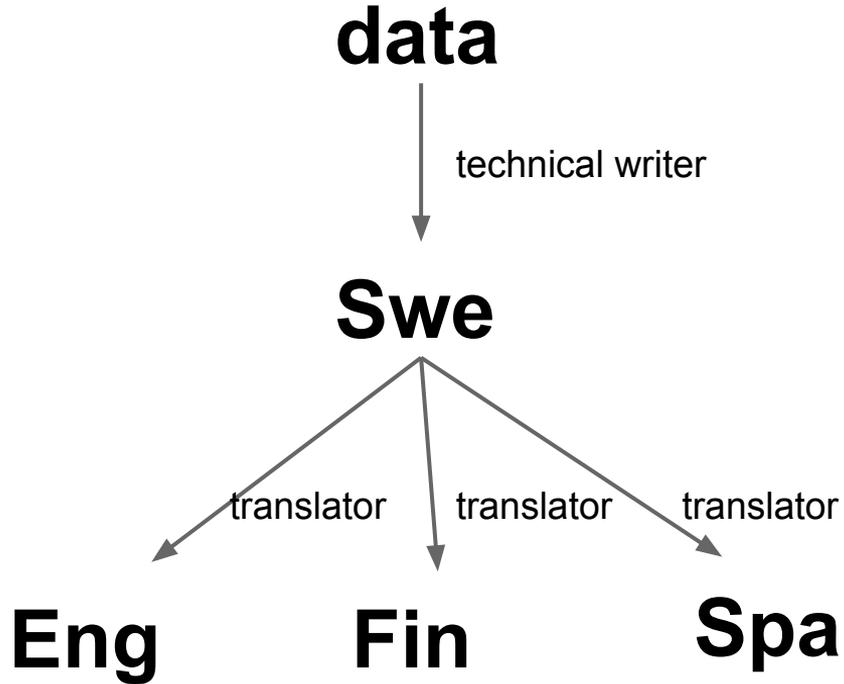
Documentation: Fin

Oven vapaa leveys on 121cm.
Kävelypinta kallistuu 0.5% siv...

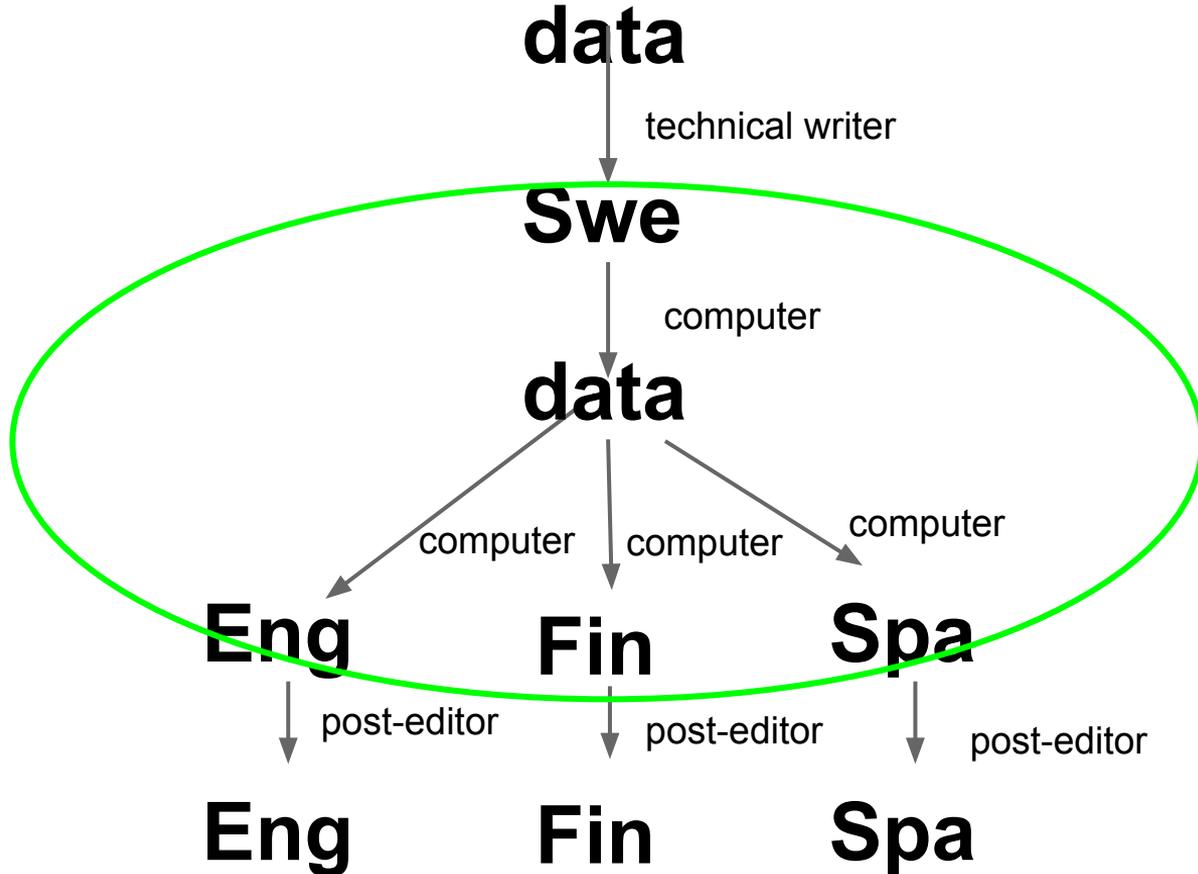
Documentation: Spa

El ancho libre de la puerta es de 121cm.
La zona peatonal se inclina 0.5% de lado

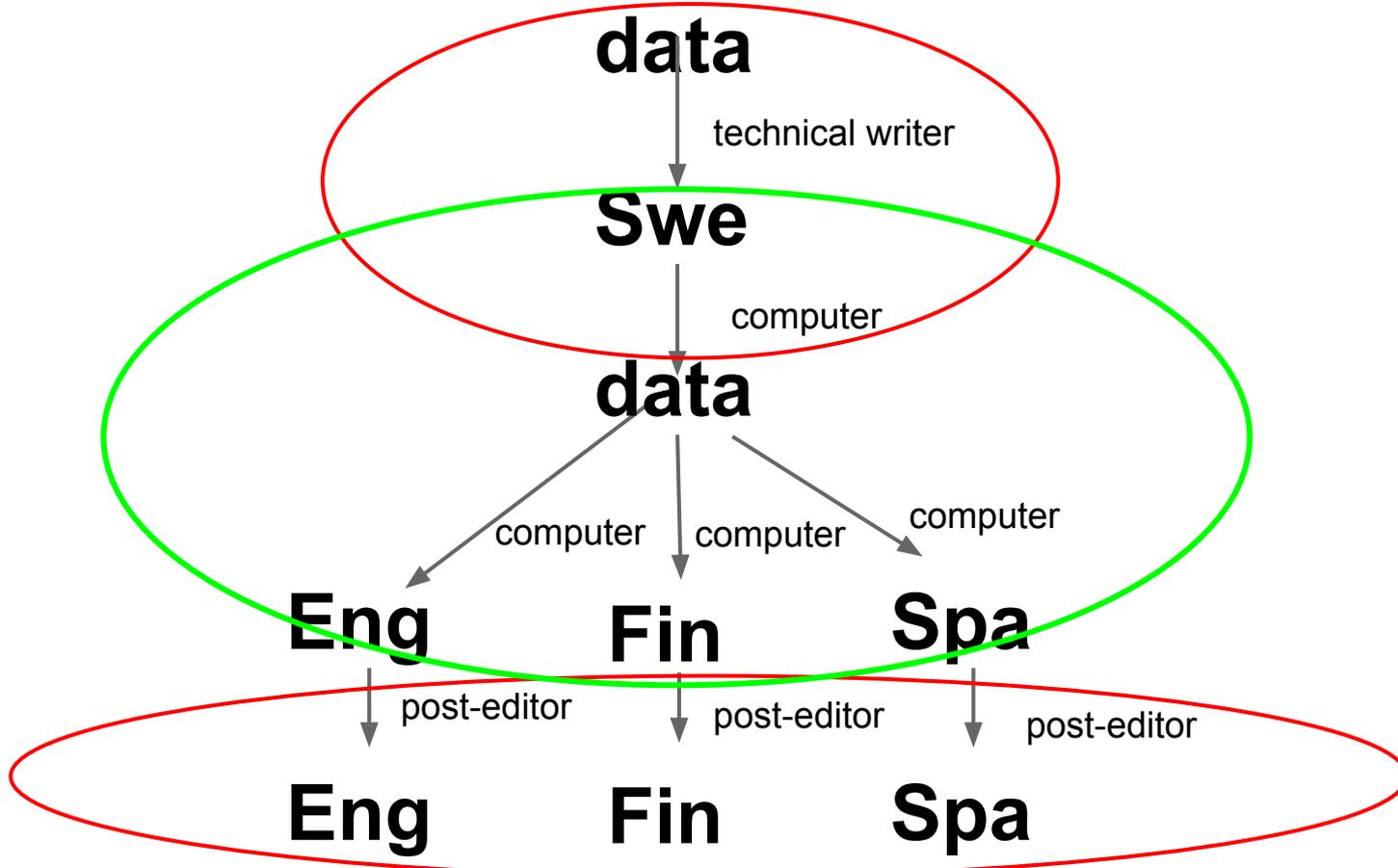
Traditional documentation



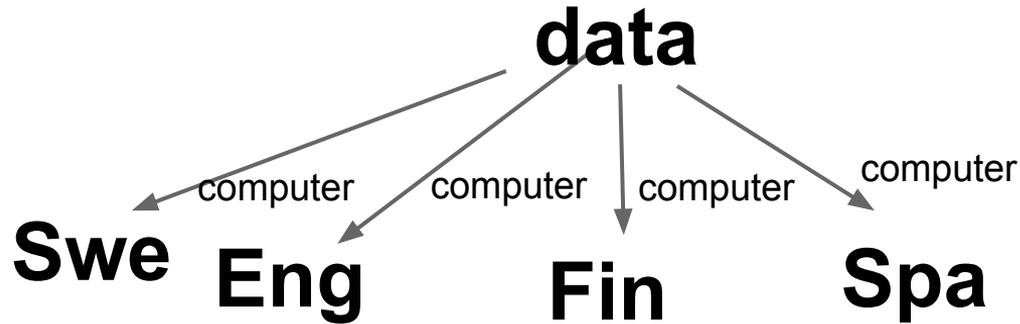
Introducing machine translation



To eliminate



Data-Driven Documentation



Advantages

Cheaper

Quicker

Better

More scalable

Cheaper

Initial cost: write the program

Later cost: mostly automatic

- post-editing at most 20% of human translation

Quicker

Translation in (almost) real time

The “almost” comes from

- new words
- post-editing need

Better

No accidental errors

Consistent terminology

More scalable

Adding new languages is easier:

- data is common to all languages

Initial effort in vocabulary

- no work with the texts themselves

How to get there

1. Extract data from texts

the door is 121cm wide

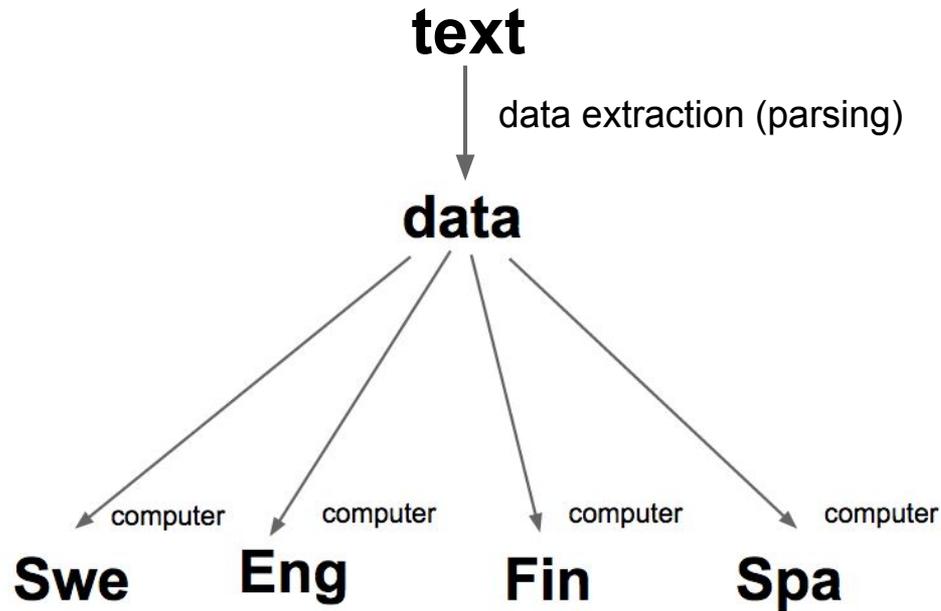
the width of the door is 121cm



door, width, 121cm

2. Support input of new information as data

Translation = Data Extraction + Data-Driven Documentation



Technology:

GF = Grammatical Framework

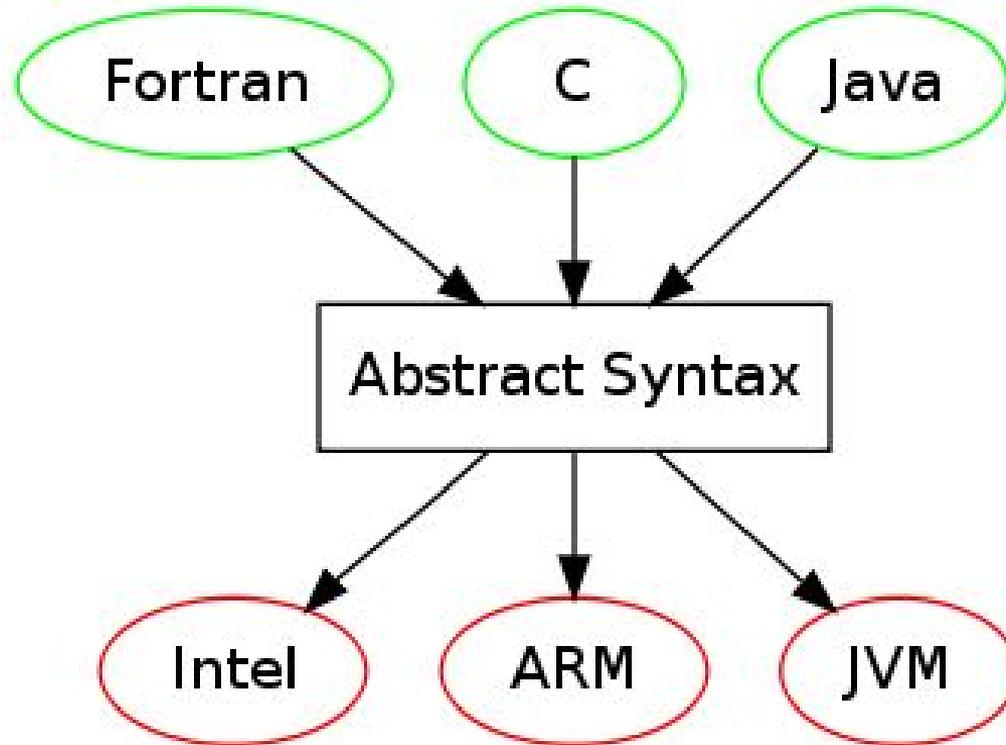
GF = Grammatical Framework

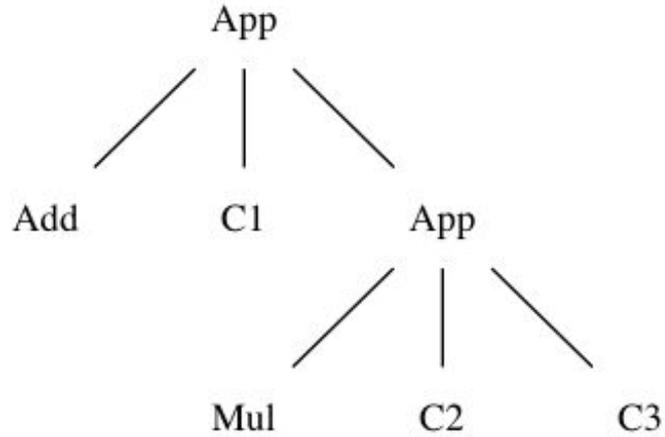
Xerox XRCE 1998, now open source

“Compiling natural language”

Library: 30 languages

Translation model: multi-source multi-target compiler



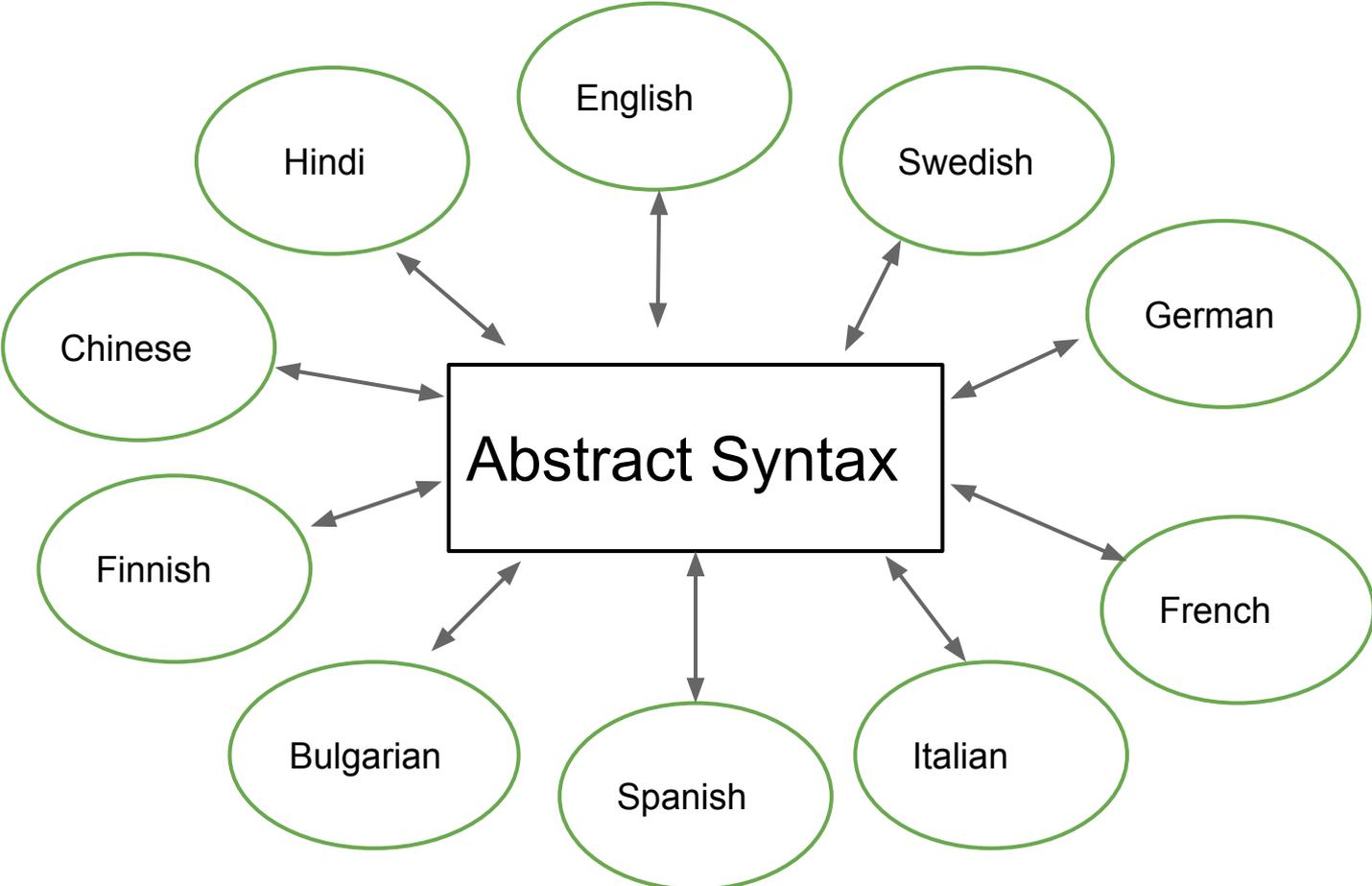


1 + 2 * 3

iconst_1
iconst_2
iconst_3
imul
iadd

(+ 1 (* 2 3))

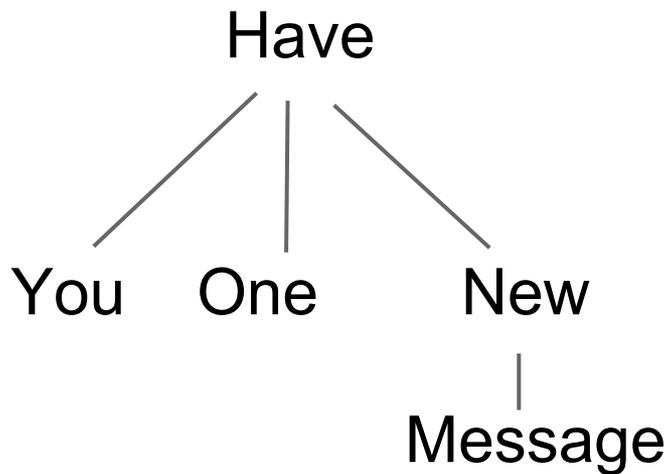
“Compiling natural language”



Abstract and concrete syntax

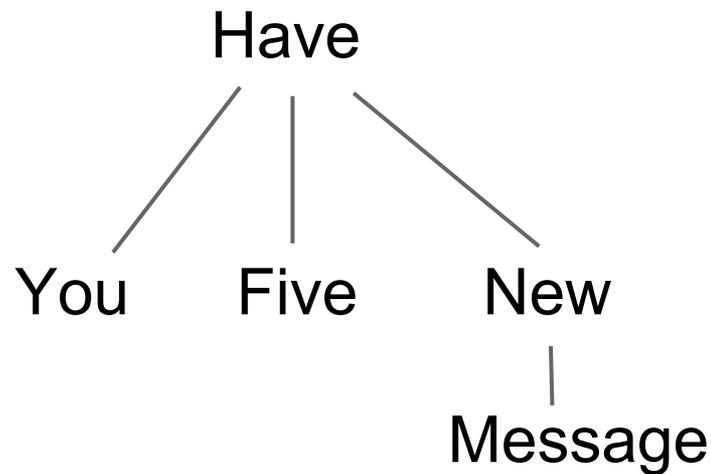
Abstract syntax: semantic structure of data

Concrete syntax: language-specific details



you have one new message

你有一个新信息



you have five new messages

你有五个新信息

What is data?

Anything that can be represented as an abstract syntax in GF!

- relational data
- Semantic Web data (OWL, RDF)
- algebraic datatypes
- logical formulas
- dependent types and lambda calculus
- Constructive Type Theory

Paintings, mathematics,...

Adam and Eve was painted by Albrecht Dürer in 1507. It measures 81 by 209 cm. This work is displayed at the Museo del Prado.

Adam and Eve a été peint par Albrecht Dürer en 1507. Il est de 81 sur 209 cm. Cette oeuvre est exposée au Musée du Prado.



Knowledge Base Results for "show everything about all paintings that are painted on canvas" (100 of many)

⚡ $\square \implies (\text{mkProp}(\text{subset}(\text{Var2Set } A)(\text{Var2Set } B))) (\text{mkProp}(\text{notprsubset}(\text{Var2Set } A)(\text{Var2Set } B)))$

⚡ ▶ ако A е подмножество на B тогава B не е грозно подмножество на D

⚡ ▶ si A és un subconjunt de B llavors B no és un subconjunt propi de D

⚡ ▶ if A is a subset of B then B is not a proper subset of D en-US

⚡ ▶ jos A on B:n osajoukko niin B ei ole D:n aito osajoukko

⚡ ▶ si A est un sous-ensemble de B alors B n'est pas un sous-ensemble propre de D

⚡ ▶ wenn A eine Teilmenge von B ist dann ist B nicht eine echte Teilmenge von D

```
PREFIX painting: <http://spraakbanken.gu.se/rdf/owl/painting.owl#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
SELECT distinct ?painting ?title ?author ?year ?length ?height ?museum
WHERE
{
  ?painting rdf:type painting:Painting ;
  rdfs:label ?title ;
  painting:hasCurrentLocation ?museum;
  painting:hasCreationDate ?date;
  painting:hasDimension ?dim ;

  painting:createdBy ?author . ?author rdfs:label ?painter .
  ?date painting:toTimePeriodValue ?year . ?dim painting:lengthValue ?length ;
  painting:heightValue ?height . ?museum rdfs:label ?loc .
```

```

TitleParagraph DefinitionTitle
DefPredParagraph type_Sort A_Var contractible_Pred (ExistCalledProp a_Var (ExpSort (VarExp A_Var)) (FunInd centre_of_contraction_Fun) (ForAllProp (BaseVar x_Var) (ExpSort (VarExp A_Var)) (ExpProp (equalExp (VarExp a_Var) (VarExp x_Var))))))
FormatParagraph EmptyLineFormat
TitleParagraph DefinitionTitle
DefPredParagraph (mapSort (mapExp (VarExp A_Var) (VarExp B_Var))) f_Var equivalence_Pred (ForAllProp (BaseVar y_Var) (ExpSort (VarExp B_Var)) (PredProp contractible_Pred (AliasInd (AppFunItnd fiber_Fun) (FunInd (ExpFun (ComprehensionExp x_Var (VarExp A_Var) (equalExp (AppExp f_Var (VarExp x_Var)) (VarExp y_Var))))))))))
DefPropParagraph (ExpProp (equivalenceExp (VarExp A_Var) (VarExp B_Var))) (ExistSortProp (equivalenceSort (mapExp (VarExp A_Var) (VarExp B_Var))))
FormatParagraph EmptyLineFormat
TitleParagraph LemmaTitle
TheoremParagraph (ForAllProp (BaseVar A_Var) type_Sort (PredProp equivalence_Pred (AliasInd (FunInd identity_map_Fun) (FunInd (ExpFun (DefExp (identityMapExp (VarExp A_Var)) (TypedExp (BaseExp (lambdaExp x_Var (VarExp A_Var) (VarExp x_Var))) (mapExp (VarExp A_Var) (VarExp A_Var))))))))))
FormatParagraph EmptyLineFormat
TitleParagraph ProofTitle
AssumptionParagraph (ConsAssumption (ForAssumption y_Var (ExpSort (VarExp A_Var)) (LetAssumption (FunInd (ExpFun (DefExp (fiberExp (VarExp y_Var) (VarExp A_Var)) (ComprehensionExp x_Var (VarExp A_Var) (equalExp (VarExp x_Var) (VarExp y_Var)))))) (AppFunItnd (fiberWrt_Fun (FunInd (ExpFun (identityMapExp (VarExp A_Var)))))) (BaseAssumption (LetExpAssumption (barExp (VarExp y_Var)) (TypedExp (BaseExp (pairExp (VarExp y_Var) (reflexivityExp (VarExp A_Var) (VarExp y_Var)))) (fiberExp (VarExp y_Var) (VarExp A_Var))))))
ConclusionParagraph (AsConclusion (ForAllProp (BaseVar y_Var) (ExpSort (VarExp A_Var)) (ExpProp (equalExp (pairExp (VarExp y_Var) (reflexivityExp (VarExp A_Var) (VarExp y_Var))) (VarExp y_Var)))) (ApplyLabelConclusion id_induction_Label (ConsInd (FunInd (ExpFun (VarExp y_Var)) (ConsInd (FunInd (ExpFun (TypedExp (BaseExp (VarExp x_Var)) (VarExp A_Var)))) (ConsInd (FunInd (ExpFun (TypedExp (BaseExp (VarExp z_Var)) (idPropExp (VarExp x_Var) (VarExp y_Var)))) BaseInd))) (DisplayExpProp (equalExp (pairExp (VarExp x_Var) (VarExp z_Var)) (VarExp y_Var))))))
ConclusionSoThatParagraph (ForConclusion (BaseVar y_Var) (ExpSort (VarExp A_Var)) (A BaseInd) (ExpProp (equalExp (VarExp u_Var) (VarExp y_Var)))) (PredProp contractible_Pri
ConclusionParagraph (PropConclusion (PredProp equivalence_Pred (FunInd (ExpFun (Type
QEDParagraph

```

Definition: A type A is contractible, if there is $a : A$, called the center of contraction, such that for all $x : A$, $a = x$.

Definition: A map $f : A \rightarrow B$ is an equivalence, if for all $y : B$, its fiber, $\{x : A \mid fx = y\}$, is contractible. We write $A \simeq B$, if there is an equivalence $A \rightarrow B$.

Lemma: For each type A , the identity map, $1_A := \lambda x:A x : A \rightarrow A$, is an equivalence.

Proof: For each $y : A$, let $\{y\}_A := \{x : A \mid x = y\}$ be its fiber with respect to 1_A and let $\bar{y} := (y, r_A y) : \{y\}_A$. As for all $y : A$, $(y, r_A y) = y$, we may apply Id-induction on y , $x : A$ and $z : (x = y)$ to get that

$$(x, z) = y$$

. Hence, for $y : A$, we may apply Σ -elimination on $u : \{y\}_A$ to get that $u = y$, so that $\{y\}_A$ is contractible. Thus, $1_A : A \rightarrow A$ is an equivalence. \square

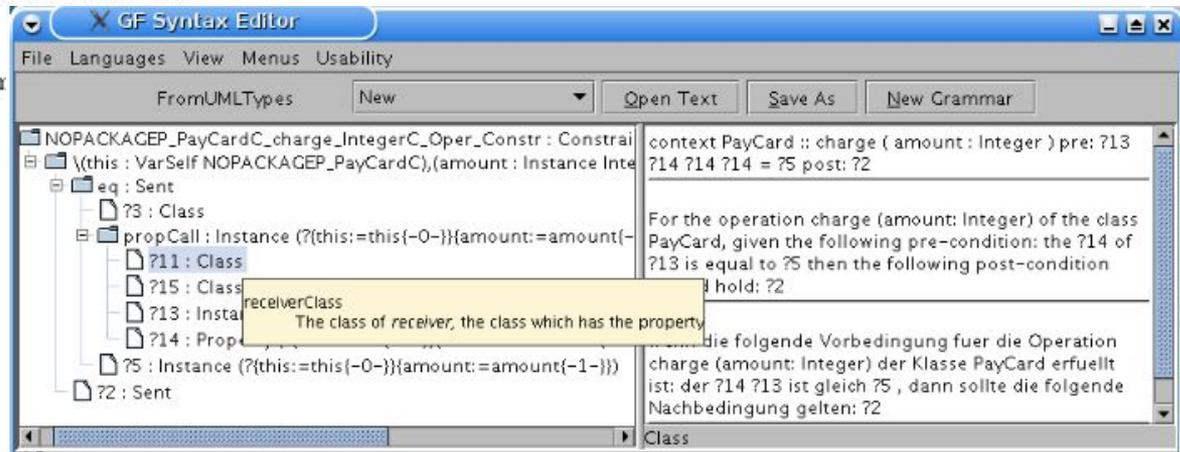
. Donc, pour les $y : A$, nous pouvons appliquer Σ -élimination sur $u : \{y\}_A$ pour obtenir que $u = y$, de façon que $\{y\}_A$ soit contractible. Alors, $1_A : A \rightarrow A$ est une équivalence. \square

GF-KeY

- if the try counter is equal to 0 then this implies that the result is equal to false
- if the following conditions are true
 - the try counter is greater than 0
 - *pin* is not equal to null
 - *offset* is at least 0
 - *length* is at least 0
 - *offset* plus *length* is at most the size of *pin*
 - the query `arrayCompare (the pin , 0 , pin , offset , length)1` on `Util` is equal to 0

then this implies that the following conditions are true

- the result is equal to true
- this owner PIN is validated
- the try counter is equal to the maximum



Some more applications

Mathematical teaching material (WebALT)

Tourist phrasebook (MOLTO)

Formal specifications (Galois)

Patent query language (Ontotext)

Museum query language and texts (Ontotext)

Business models (Be Informed)

Medical examination journals (Lingsoft)

Speech commands in cars (Talkamatic)

Accessibility database (Digital Grammars/TD)

Norwegian

Danish

Afrikaans

English Swedish German Dutch

French Italian Spanish Catalan

Bulgarian Finnish Estonian

Japanese Thai Chinese Hindi

Latvian Mongolian Urdu Punjabi Sindhi

Greek Maltese Nepali Persian

Latin Turkish

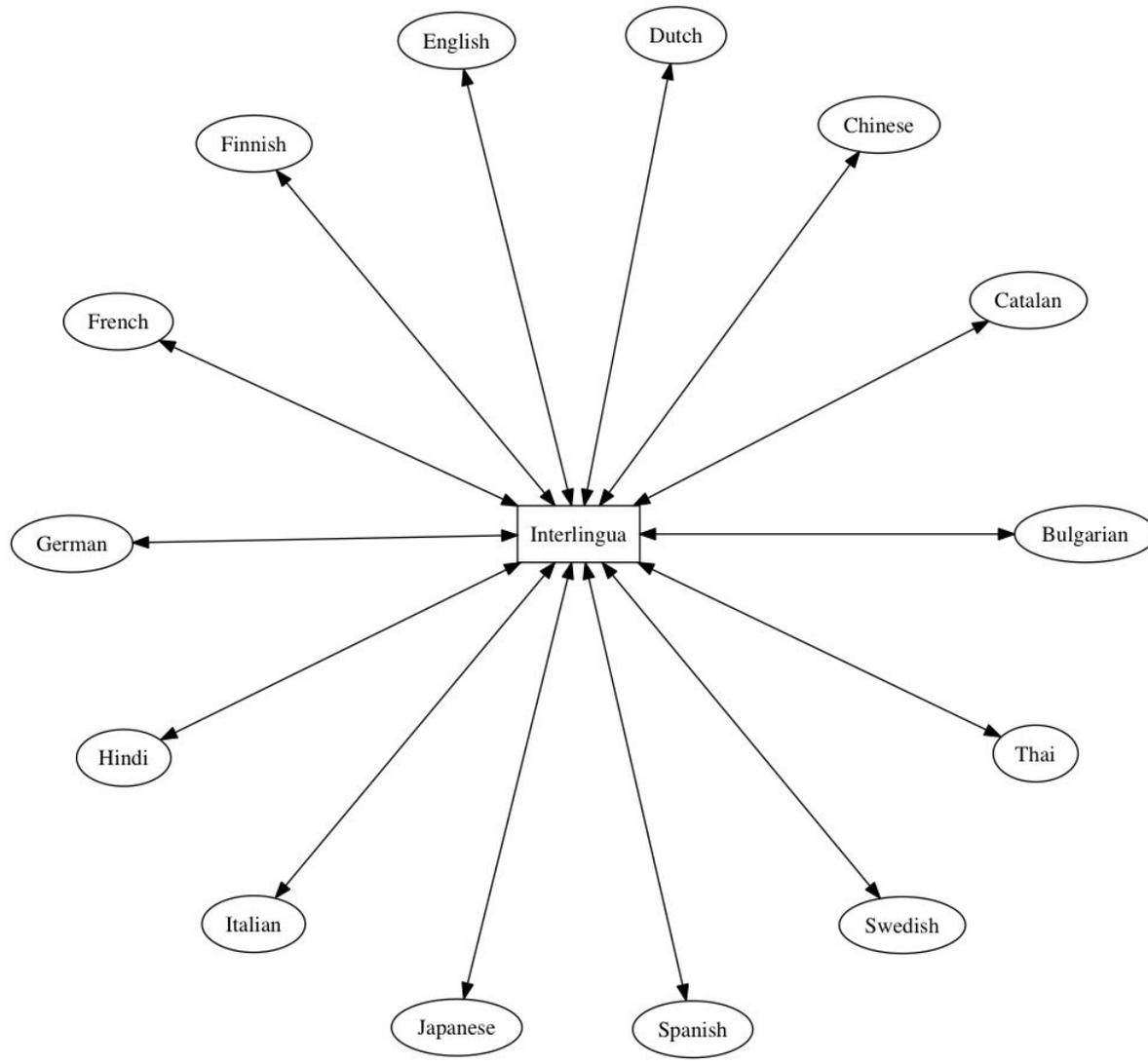
Hebrew Arabic Amharic

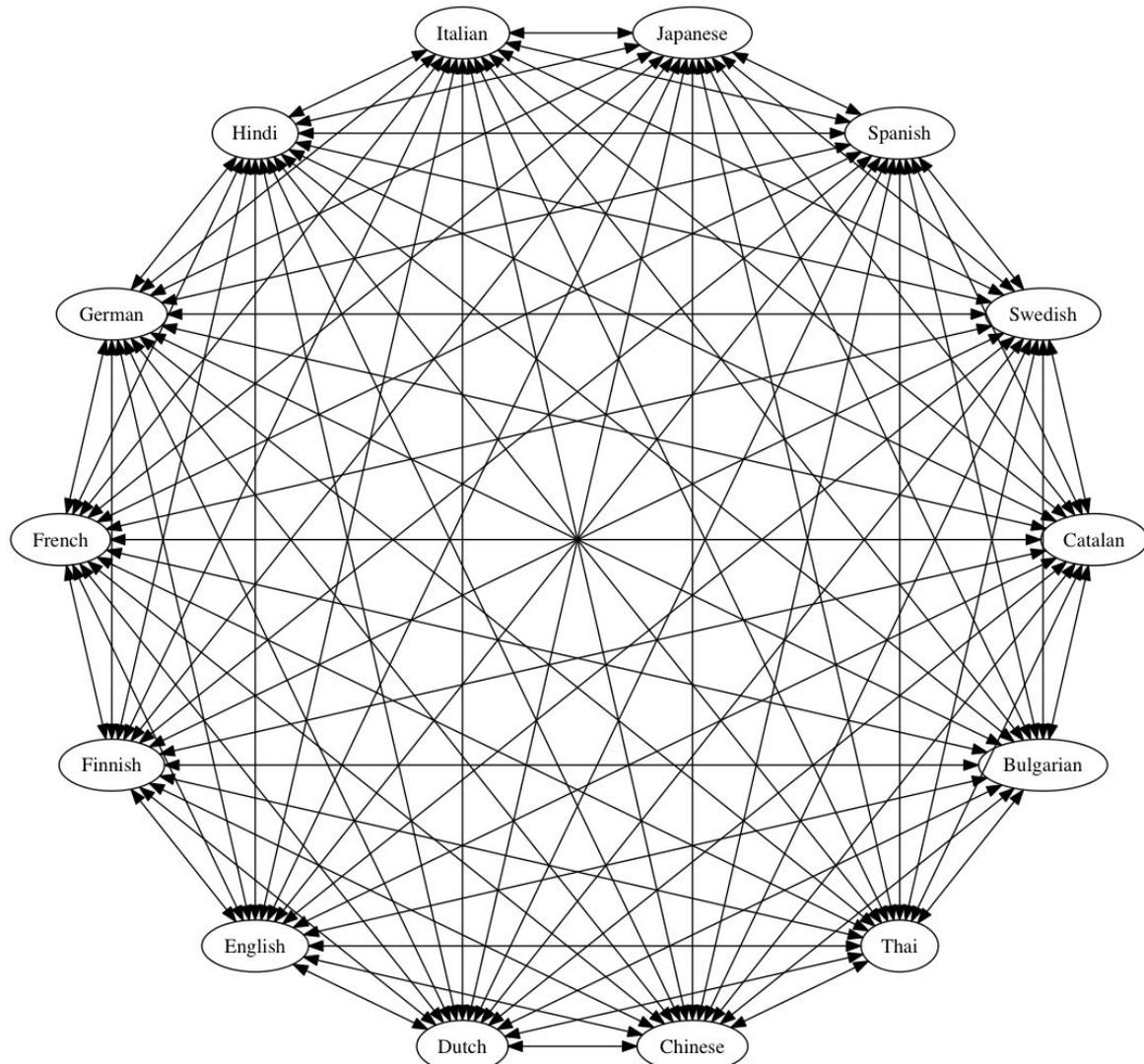
Swahili

Romanian

Polish

Russian

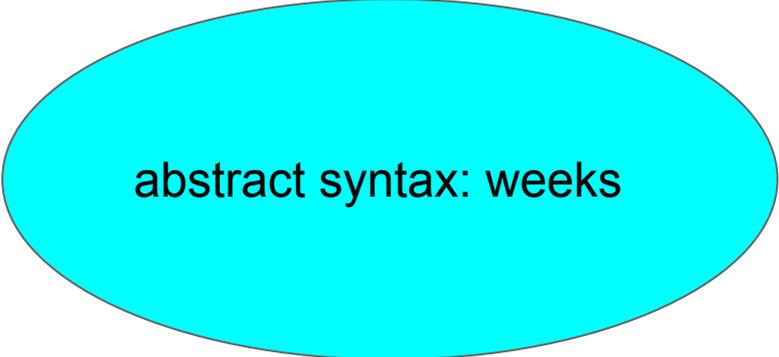




Domain adaptation

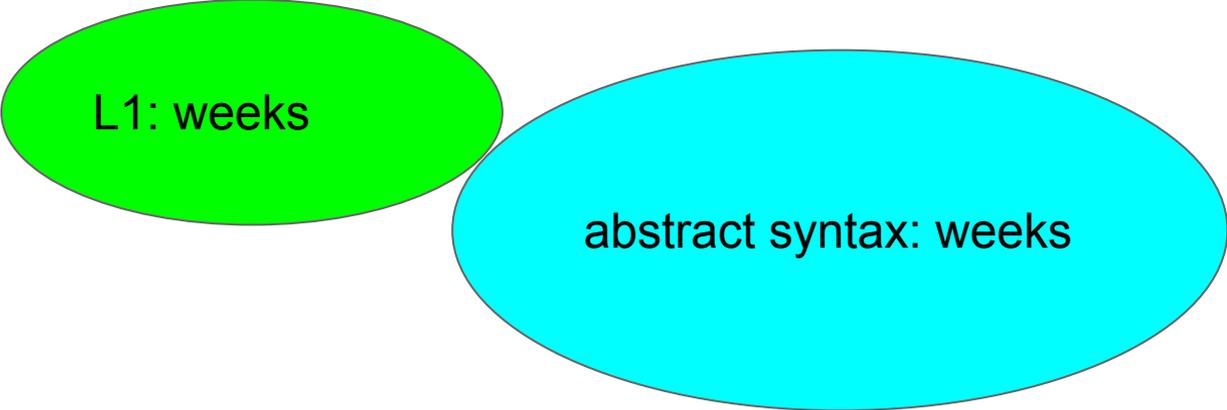
1. Build an abstract syntax to model the domain.
 - The biggest one-time cost.
2. Build concrete syntaxes for the languages you want to cover.
 - Cost goes down as languages are added.

Building effort



abstract syntax: weeks

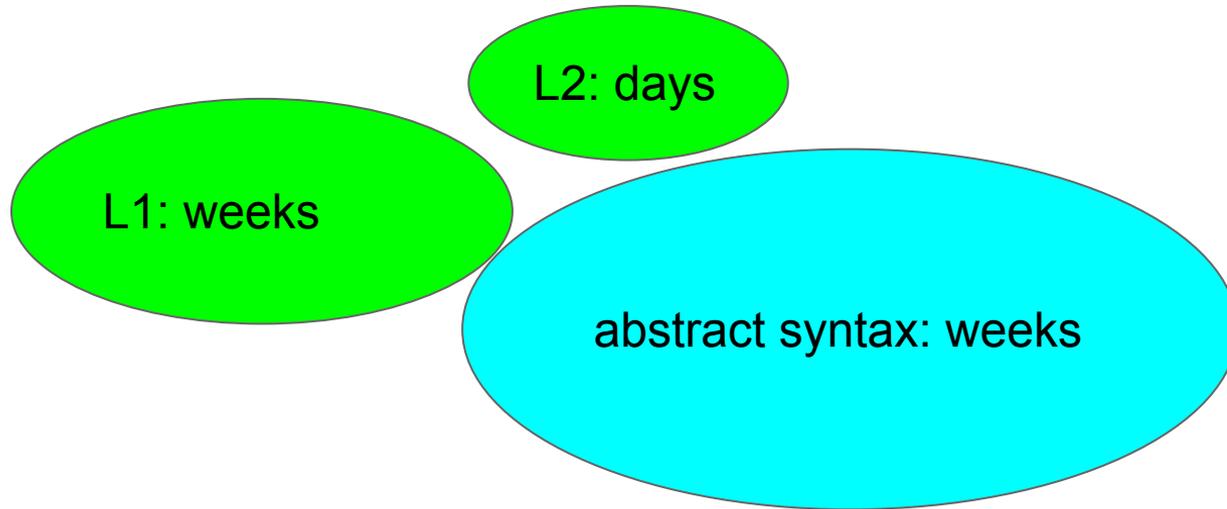
Building effort



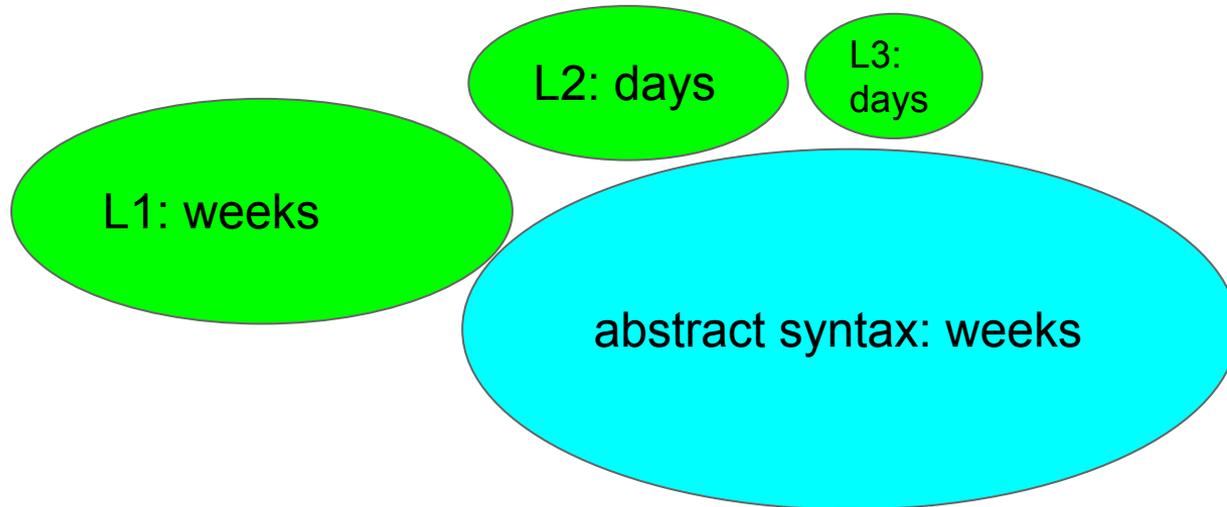
L1: weeks

abstract syntax: weeks

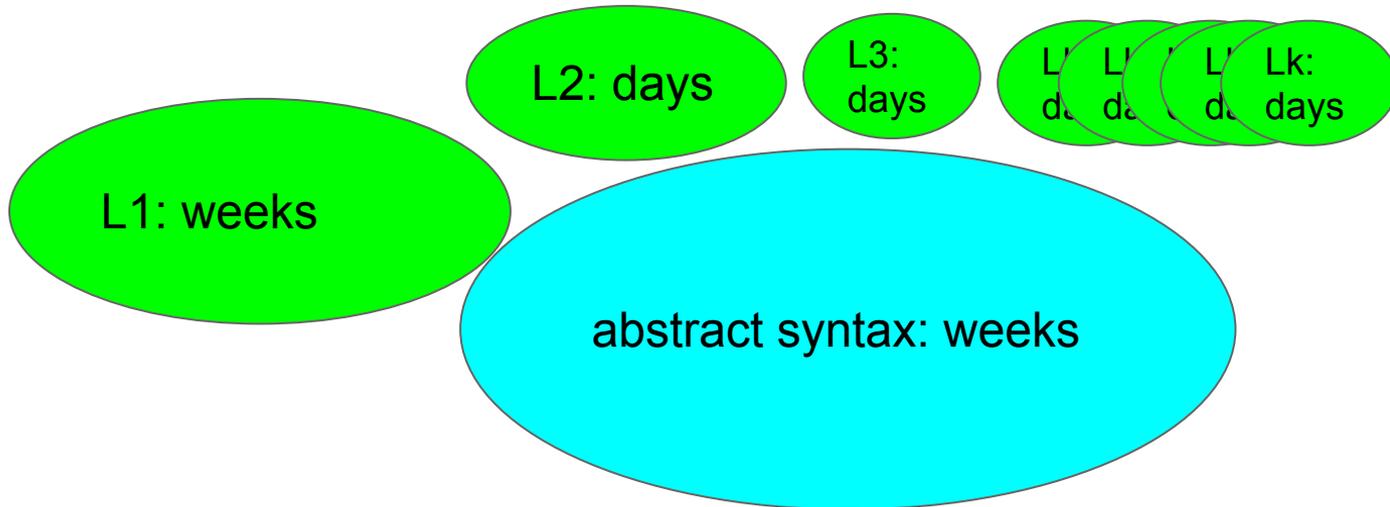
Building effort



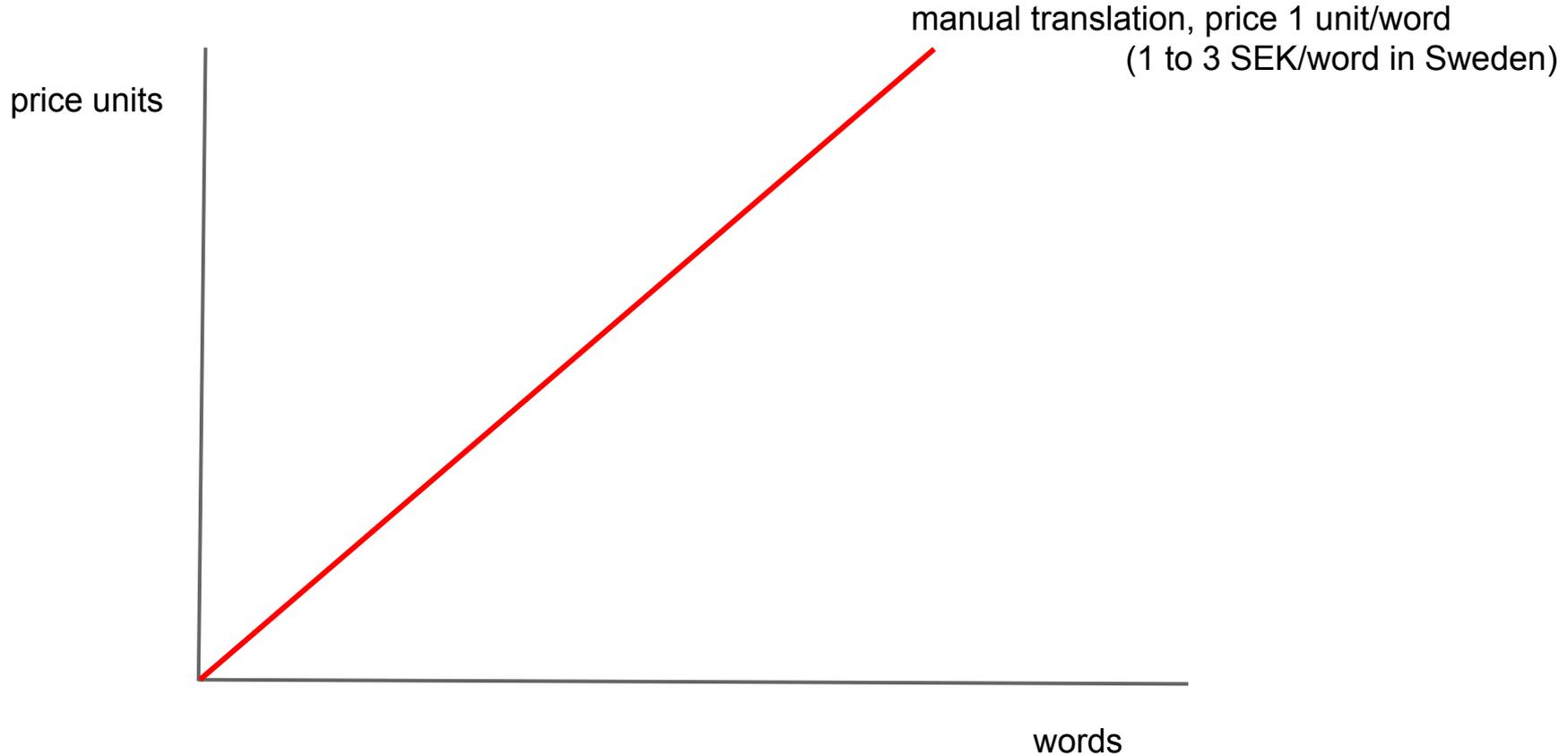
Building effort



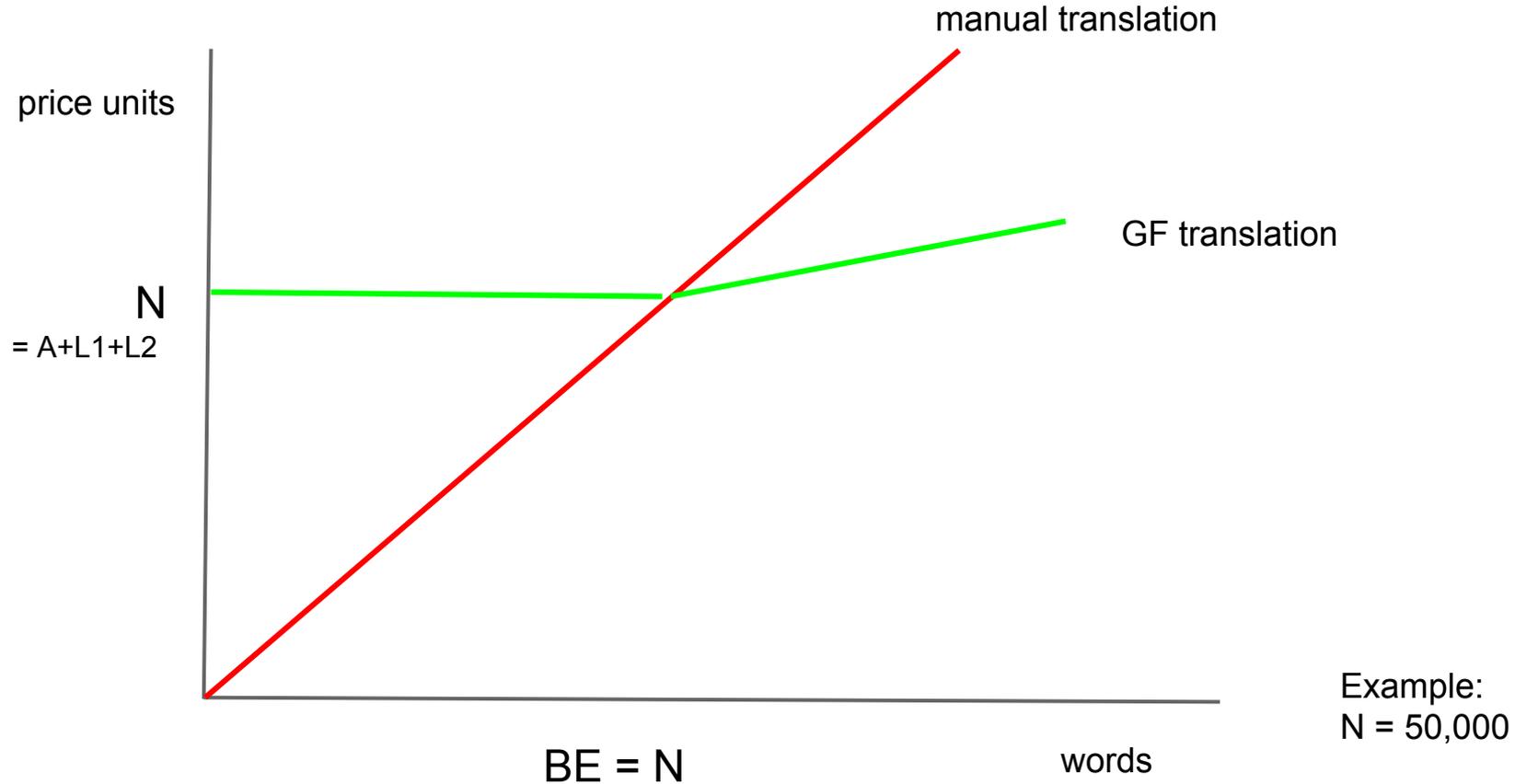
Building effort



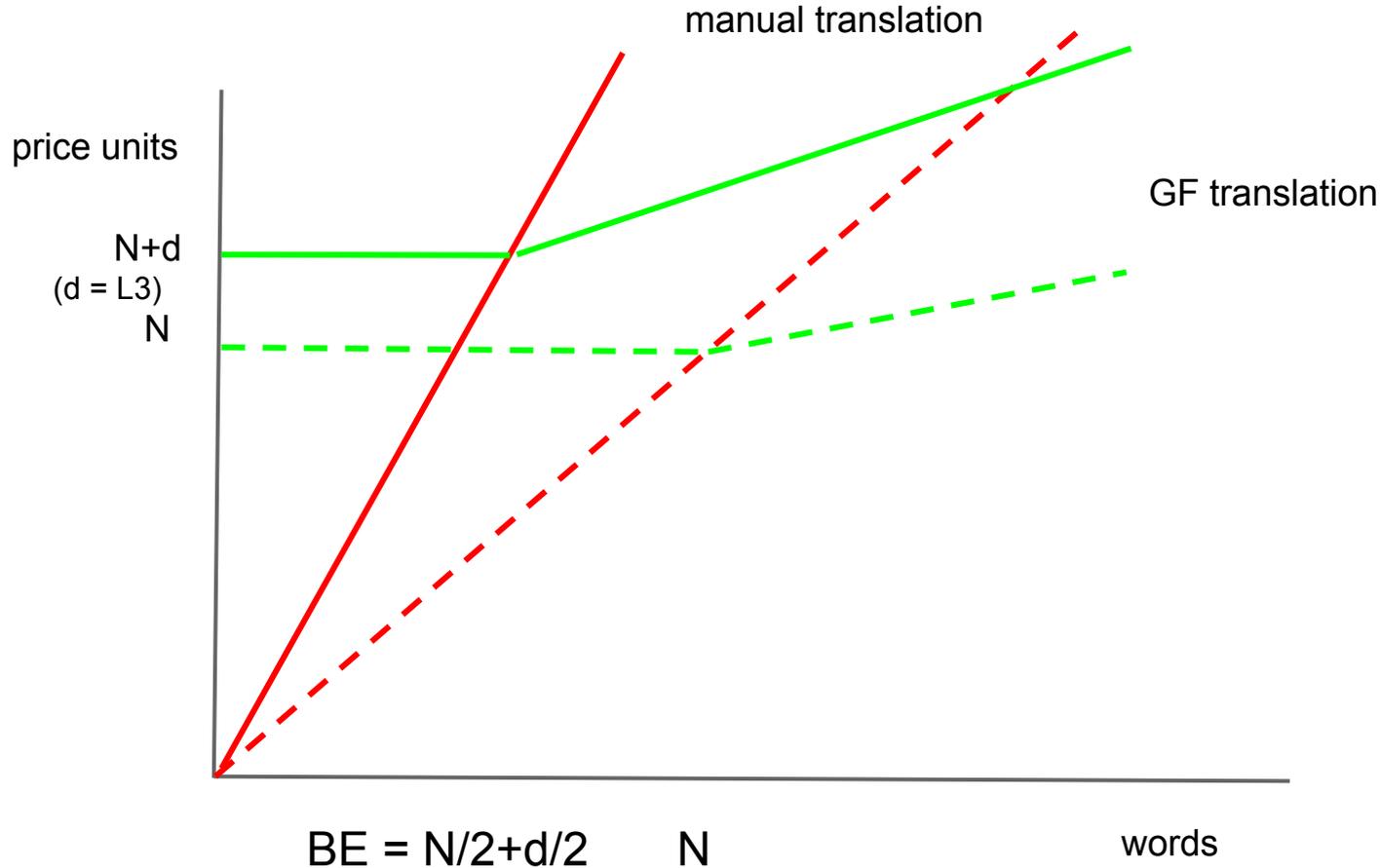
Price of translation, 1 target language



Break-even point, 1 target language

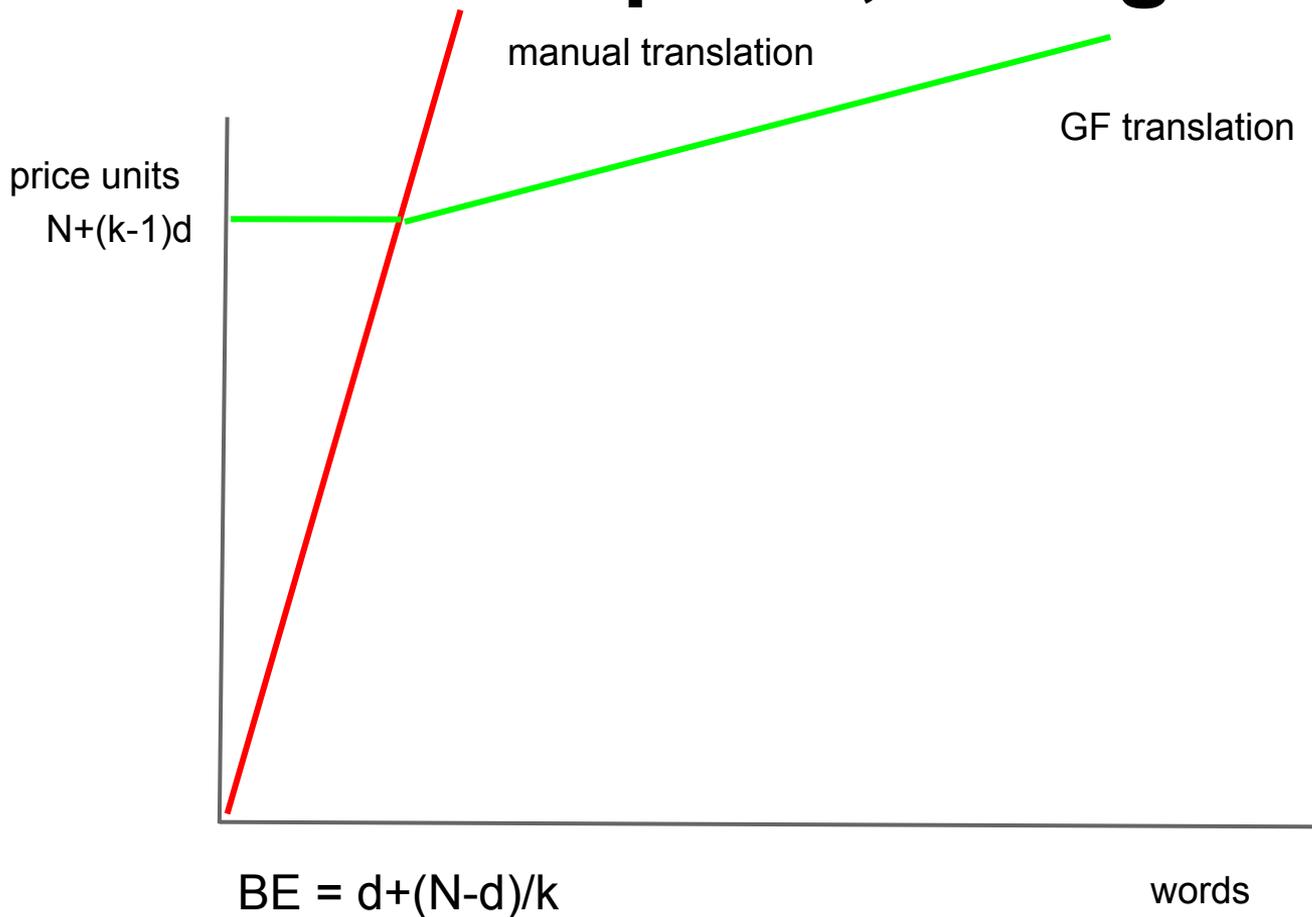


Break-even point, 2 target languages



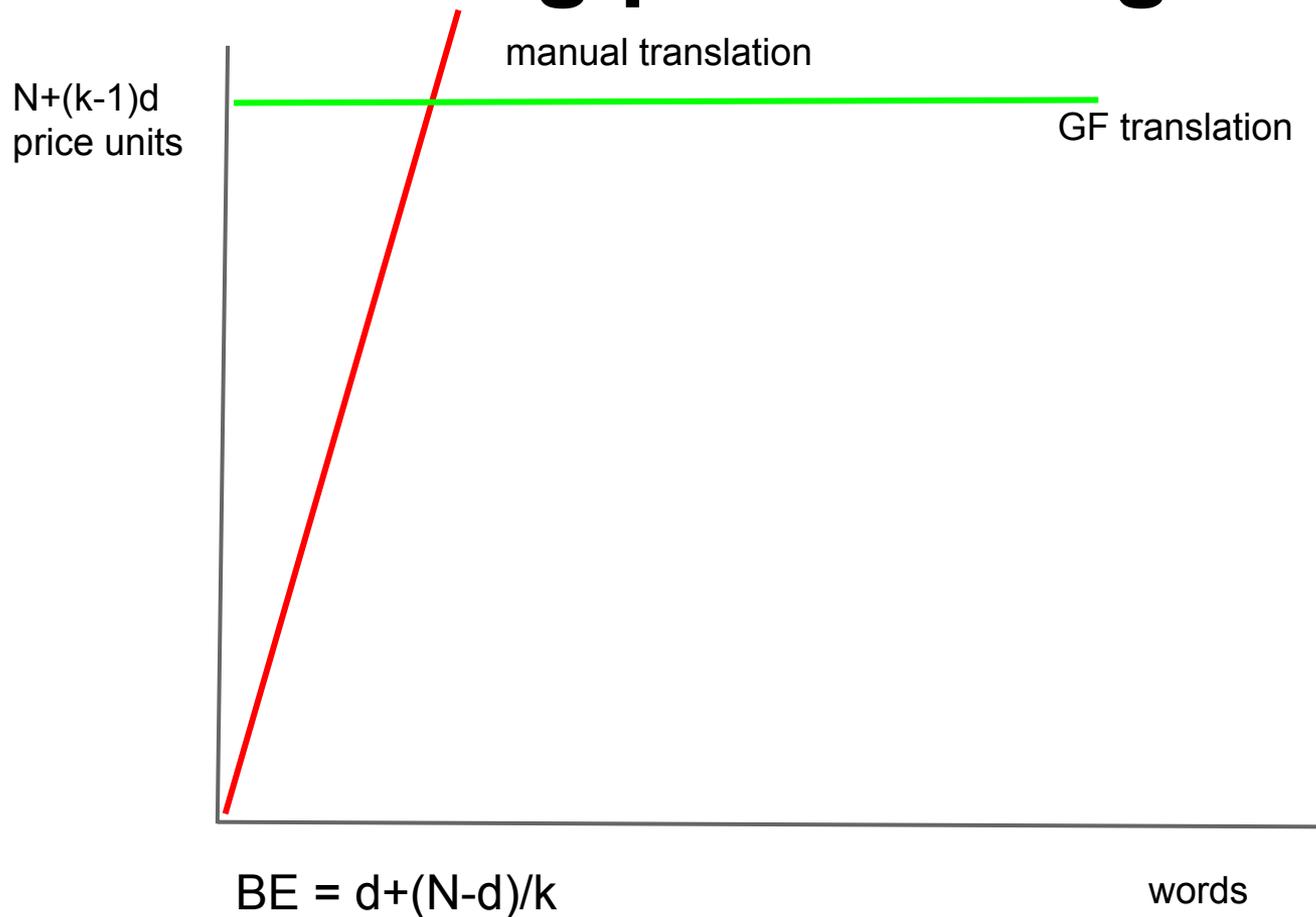
Example:
 $N = 50,000$, $d = 20,000$
 $BE = 35,000$

Break-even point, k target languages



Example: $k=10$
 $N = 50,000$, $d = 20,000$
 $BE = 23,000$

Eliminating post-editing



Example: $k=10$
 $N = 100,000$, $d = 30,000$
 $BE = 37,000$

The biggest challenge

To reach the customers who need this...

The biggest challenge

To reach the customers who need this...

... or at least something like this

We need to translate to many languages

We need to update the content often

Google translate is too low quality

Human translation is too expensive or too slow

Candidates

Health care

E-commerce

Legal documentation, contracts

Technical documentation, manuals

Robot journalism

Existing customers



- Svängrumsytan utanför dörren lutar 1% i sidled.
- The turning space outside the gate tilts 1% sideways.
- Kääntymätila oven ulkopuolella kallistuu 1% sivusuunnassa.
- Der Schwenkbereich außerhalb der Tür neigt sich um 1% seitlich

UttSTD (PredUttTD (AdvNPTD (DetCNTD (DetQuant DefArt NumSg) (UseNTD svängrumsyta_NTD)) (PrepNPTD utanför_Prep (DetCNTD (DetQuant DefArt NumSg) (UseNTD dörr_NTD)))) (AdvVPTD (luta_VPTD (procentMeasure 1)) i_sidled_AdvTD))

Talkamatic
FREE DIALOGUE



digital **G**rammars
Language technology to rely on.

`next_membership_level_sys_answer silver (next_membership_points_sys_answer integer0_99_50)`

`test_mockup_travelChi: 您有五十个常旅客点符合会员条件, 您现在是在伦敦.`

`test_mockup_travelDut: je hebt vijftig punten nodig om het zilveren niveau te bereiken`

`test_mockup_travelEng: you need fifty points to reach silver level`

`test_mockup_travelFin: sinä tarvitset viisikymmentä pistettä päästäksesi hopeatasolle`

`test_mockup_travelFre: tu as besoin de cinquante points pour atteindre le niveau argent`

`test_mockup_travelGer: Sie brauchen fünfzig Punkte um das Silberniveau zu erreichen`

`test_mockup_travelIta: avete bisogno di cinquanta punti per raggiungere il livello argento`

`test_mockup_travelSpa: necesitas cincuenta puntos para llegar al nivel plata`

Home Visual Scenes

Brian - Visual Scenes





Välj vad du vill prata om 


Huvudmeny


Skrivet språk


Talat språk


Logo mat


Trädgården


Fika


Se på TV

I don't know where my wife is	non riesco a trovare mia moglie	jag vet inte var min fru är	ich weiss nicht, wo meine Frau ist	Ne vem, kje je moja žena
I don't know where my husband is	non riesco a trovare mio marito	jag vet inte var min man är	ich weiss nicht ,wo mein Mann ist	Ne vem, kje je moj muž
I'm hurt	mi sono ferito	jag har skadat mig	ich bin verletzt	Poškodovan sem
I need a doctor	ho bisogno di un medico	jag behöver en läkare	ich brauche einen Arzt	rabim zdravnika
I'm allergic to penicillin	Sono allergico alla penicillina	jag är allergisk mot penicillin	ich bin allergisch gegen Penicilin	sem alergičen na penicilin
I'm dizzy	mi gira la testa	jag är yr	mir ist schwindelig	Sem omotična
I have nausea	ho la nausea	jag mår illa	ich habe Übelkeit	Čutim slabost, počutim se slabo
my arm hurts	mi fa male il braccio	jag har ont i armen	mein Arm tut weh	me roka boli

Data-Driven Question Answering

A derived product

I want to go from
Pudong Airport to
Hongqiao Station.

I want to go from
Pudong Airport to
Hongqiao Station.

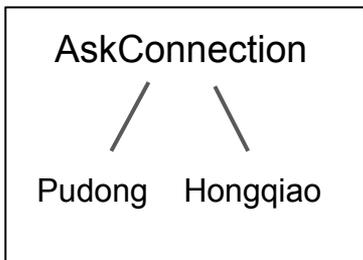
parsing

AskConnection

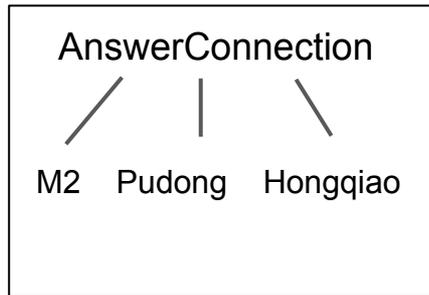
Pudong Hongqiao

I want to go from
Pudong Airport to
Hongqiao Station.

parsing

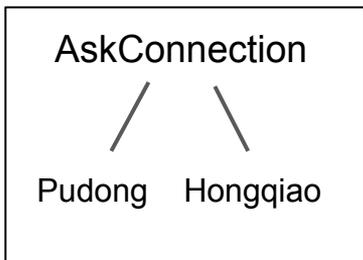


query engine

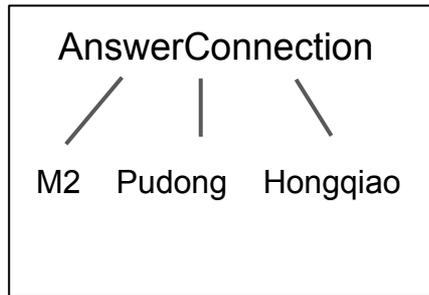


I want to go from
Pudong Airport to
Hongqiao Station.

parsing



query engine

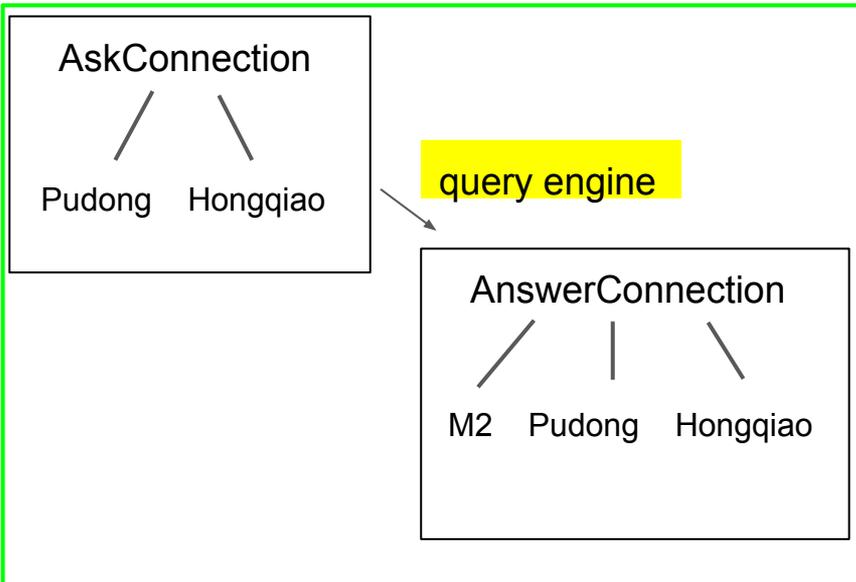


linearization

Take Metro line 2
from Pudong Airport
to Hongqiao Station.

I want to go from
Pudong Airport to
Hongqiao Station.

parsing



query engine

linearization

Take Metro line 2
from Pudong Airport
to Hongqiao Station.

从浦东机场到虹桥站怎么走？

parsing

AskConnection

Pudong Hongqiao

query engine

AnswerConnection

M2 Pudong Hongqiao

linearization

在浦东坐2号地铁到虹桥站

Kuinka pääsee
Pudongin lentokentältä
Hongqiao-asemalle?

parsing

AskConnection

Pudong Hongqiao

query engine

AnswerConnection

M2 Pudong Hongqiao

linearization

Mene metrolla 2
Pudongin
lentokentältä
Hongqiao-asemalle.